



Soumya Mukherjee

CH24 M571

ML

ID 5001 W

DATE 11/4/25

PLACE Kolkata

VOP Assignment - 2

(1) we have 2 classes $y_i = 1$ & $y_i = -1$

Given 4 points in 2D feature space

$$x_1 = (2, 2), y_1 = +1$$

$$x_2 = (4, 4), y_2 = +1$$

$$x_3 = (1, -1), y_3 = -1$$

$$x_4 = (-3, -3), y_4 = -1$$

Given normal vector to hyperplane at
 $0 \Rightarrow w_0 = [0, 0]$.Bias at 0 $\Rightarrow b_0 = 0$.

Update rule is

$$w^{t+1} = w^t + y_i x_i$$

$$b^{t+1} = b^t + y_i$$

When prediction $\neq y_i$, we need to
update weight 'w' & bias 'b'.

We know, Perceptron update rule is

If $w^T x_i + b \geq 0$, predict 1

else if $w^T x_i + b < 0$, predict -1

Epoch - 1

(a) $x_1 = (2, 2)$.

$$(0,0) + (0,0) \cdot (2,2) + 0 = d + w \cdot x$$

$$w \cdot x + b = (0,0) \cdot (2,2) + 0$$

$= 0 + 0 \rightarrow \text{Predict: } +1$
Given -1 misclassified as +1 (correct Prediction)

(b) $x_2 = (4, 4)$

$$w \cdot x + b = (0,0) \cdot (4,4) + 0$$

$= 0 + 0 = 0 \rightarrow \text{Predict: } +1$
Given +1 as -1 misclassified as +1 correct Prediction.

(c) $x_3 = (1, -1)$

$$w \cdot x + b = (0,0) \cdot (1, -1) + 0$$

$= 0 \rightarrow \text{Predict: } +1$

Given -1 wrong prediction

(d) $x_4 = (-3, -3)$

$$w \cdot x + b = (0,0) \cdot (-3, -3) + 0$$

$= 0 \rightarrow$



DATE _____

PLACE _____

Since, wrong prediction, we would want to update weights & bias.

$$\bullet \quad w = (0, 0) + (1, -1) \cdot (-1) = (-1, 1)$$

$$b = 0 + (-1) = -1$$

(d) $x_4 = (-3, -3)$

$$wx + b = (-1, 1) \cdot (-3, -3) + (-1)$$

$$= 3 + (-3) - 1$$

$$= -1 \quad \text{Prediction} = -1 \quad \text{since } < 0.$$

: Correct Prediction

Epoch - 2

(a) $x_1 = (2, 2), y = +1 \quad w = (-1, 1) \quad b = -1$

$$wx + b = (-1, 1) \cdot (2, 2) + -1$$

$$= -1$$

wrong Prediction



DATE _____

PLACE _____

Performing update operation..

$$\text{new } w = (-1, 1) + (+1)(2, 2)$$

$$= (1, 3)$$

$$b = -1 + (+1) = 0$$

(b) $x_2 = (4, 4), y = +1$

$$wx + b = (1, 3) \cdot (4, 4) + 0$$

$\text{Ans: } 16 > 0 \Rightarrow \text{Predict: 1}$

∴ Thus, prediction is correct.

(c) $x_3 = (1, -1), y = -1$

$$\Rightarrow wx + b = (1, 3) \cdot (1, -1) + 0$$

$$= 1 - 3 + 0 = -2 < 0.$$

Predict -1, Prediction: Correct

(d) $x_4 = (-3, -3), y = -1$

$$wx + b = (1, 3) \cdot (-3, -3) + 0$$

$$= -3 - 9 = -12 < 0,$$

Predict: correct (-1)

Thus after 2 updates,

$$w = (1, 3) \quad \& \quad b = 0$$

Ans (a)



DATE _____

PLACE _____

(b) The data is linearly separable if

$$y(w^T x + b) > 0 \text{ for all points}$$

i.e. the predicted signs matches the actual class label y for all data points.

using $w = (1, 3)$ & $b = 0$, if we are able to correctly predict all the 4 cases, we can claim it to be linearly separable.

Point	Given label	$w \cdot x + b$	Prediction	Correct
(2, 2)	+1	$1 \times 2 + 3 \times 2 = 8 > 0$	+1	Yes
(4, 4)	+1	$4 \times 1 + 3 \times 4 = 16 > 0$	+1	Yes
(1, -1)	-1	$1 \times 1 + 3 \times (-1) = -2 < 0$	-1	Yes
(-3, -3)	-1	$1 \times (-3) + 3 \times (-3) = -12 < 0$	-1	Yes



DATE _____

PLACE _____

Since, all predictions match the labels, the data is linearly separable by the hyperplane $w^T x_i + b = 0$.

(c) we need to calculate γ -margin for each points.

$$\boxed{\gamma_i = \frac{|y_i(w^T x_i + b)|}{\|w\|}}$$

$$\text{with given } w = (1, 3) \rightarrow \|w\| = \sqrt{1+9} = \sqrt{10}$$

$$b = 0$$

$$(a) x = (2, 2), y = +1$$

$$\gamma_i = \left(+1 \times (1, 3)^T \cdot (2, 2) + 0 \right) / \sqrt{10}$$

$$= 8 / \sqrt{10}$$

$$(b) x = (4, 4), y = +1$$

$$\gamma_i = \left(+1 \times (1, 3)^T \cdot (4, 4) + 0 \right) / \sqrt{10}$$

$$= 16 / \sqrt{10}$$

$$(c) x = (1, -1), y = -1$$

$$\gamma_i = \left(-1 \times (1, 3)^T \cdot (1, -1) + 0 \right) / \sqrt{10}$$

$$= -2 / \sqrt{10}$$



DATE _____

PLACE _____

$$(a) \quad x_0 = (-3, -3), \quad y = -1$$

$$\gamma_i = -1 \times ((1, 3)^T \cdot (-3, -3)^T + 0) / \sqrt{10}$$

$$= 12 / \sqrt{10}$$

Thus, among the 4 points given,

(c) $x = (1, -1)$ has the lowest margin.

(Ans) —



DATE _____

PLACE _____

Question - 2

SVM

Decision boundary

Margin and MVE

Optimal weight vector

Margin

MVE

Optimal weight vector

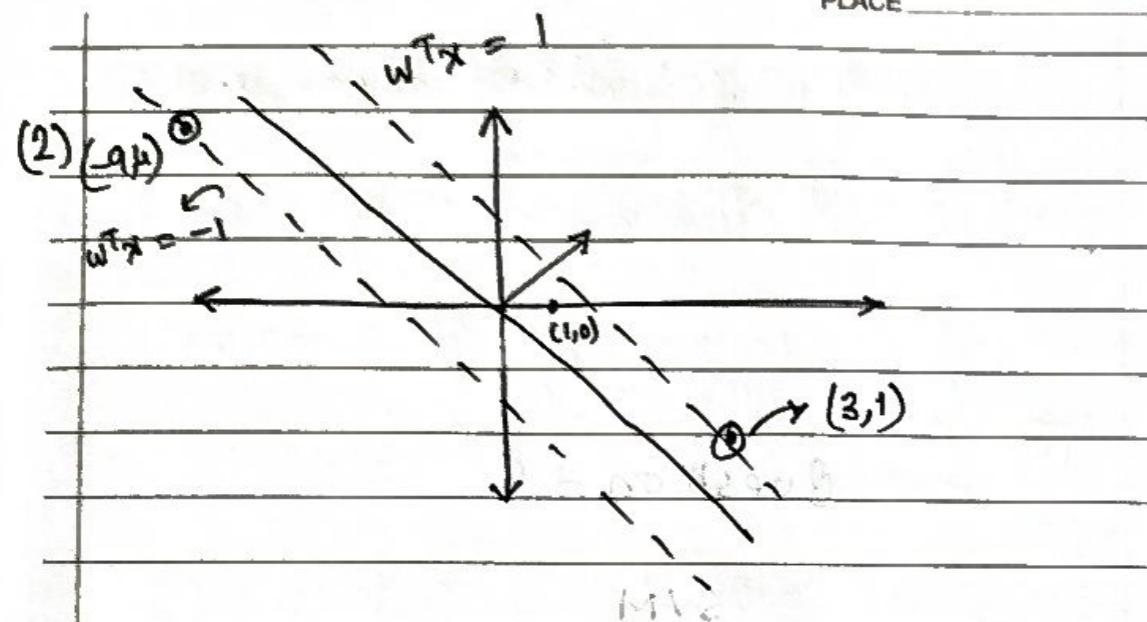
MVE

Optimal weight vector

MVE

Optimal weight vector

MVE



(1,0) is a test point.

Given hard-margin SVM i.e perfectly separates the data.

(a) Supporting hyperplanes are

$$w^T x = 1 \text{ for class } +1 \quad \& \quad w^T x = -1 \\ (\text{green}) \qquad \qquad \qquad (\text{red})$$

The decision boundary is $w^T x = 0$.
which is exactly half way between the
two supporting hyperplanes.

(b) width of separation between 2 hyperplanes is given by,

$$\frac{2}{\|w\|}$$

Given, $(-9,4)$ is on line $w^T x = -1$
 $\& (3,1)$ is on line $w^T x = 1$.

$$\text{if } w = [w_1, w_2]$$

then,

$$[w_1, w_2]^T [-9, 4] = -1$$

$$\& [w_1, w_2]^T [3, 1] = 1.$$

$$\text{i.e., } -9w_1 + 4w_2 = -1 \rightarrow ①$$

$$\& 3w_1 + w_2 = 1 \rightarrow ②$$

Multiplying ② by 3 & summing to 1.

$$9w_1 + 3w_2 = 3$$

$$+ \quad \quad \quad 7w_2 = 2$$

$$\therefore w_2 = 2/7 = 6/21$$

$$\text{thus, } 3w_1 = 1 - 2/7 = 5/7$$

$$w_1 = \frac{5}{21}$$



DATE _____

PLACE _____

$$\text{Thus, } w = \begin{bmatrix} 5/21 & 6/21 \end{bmatrix}$$

Calculating L2 norm $\|w\|_2$

$$\begin{aligned}\|w\| &= \sqrt{\left(\frac{5}{21}\right)^2 + \left(\frac{6}{21}\right)^2} = \frac{1}{21} \sqrt{25+36} \\ &= \frac{1}{21} \sqrt{61}.\end{aligned}$$

Thus, width of separation

$$= \frac{2}{\|w\|} = \frac{42}{\sqrt{61}} \quad (\text{Ans})$$

→ Equation of decision boundary.

$$w_1 x_1 + w_2 x_2 = 0$$

$$\text{or } \frac{5}{21} x_1 + \frac{6}{21} x_2 = 0$$

$$\text{or, } 5x_1 + 6x_2 = 0 \quad (\text{Ans})$$



DATE _____

PLACE _____

k-means

(3)

Initial points

$$A_1 = (2, 10) \rightarrow C_1$$

$$A_2 = (2, 5)$$

$$A_3 = (8, 4)$$

$$A_4 = (5, 8) \rightarrow C_2$$

$$A_5 = (7, 5)$$

$$A_6 = (6, 4)$$

$$A_7 = (1, 2) \rightarrow C_3$$

$$A_8 = (4, 9).$$

Initial clusters

$$= A_1, A_4, A_7$$

Calculating distance to each point

pt	Dist to A1	Dist to A4	Dist to A7	Assign
A1	0	$\sqrt{3} = 3.61$	$\sqrt{61}$	C1
A2	$\sqrt{25} = 5$	$\sqrt{18} = 4.24$	$\sqrt{16}$	C3
A3	$\sqrt{36} = 6$	$\sqrt{25} = 5$	$\sqrt{53} = 7.28$	C2
A4	$\sqrt{13} = 3.61$	0	$\sqrt{52} = 7.21$	C2
A5	$\sqrt{50} = 7.07$	$\sqrt{13} = 3.61$	$\sqrt{45} = 6.71$	C2
A6	$\sqrt{52} = 7.21$	$\sqrt{17} = 4.12$	$\sqrt{59} = 5.38$	C2
A7	$\sqrt{65} = 8.06$	$\sqrt{52} = 7.21$	0	C3
A8	$\sqrt{5} = 2.24$	$\sqrt{2} = 1.41$	$\sqrt{58} = 7.62$	C2

May

June



DATE

PLACE

After 1 epoch, new clusters become.

C₁ : A₁

C₂ : A₃, A₄, A₅, A₆, A₈.

C₃ : A₂, A₇.

New cluster centres:

$$C_2: \frac{x: 8+5+7+6+4}{5} = \frac{30}{5} = 6.$$

$$Y: \frac{4+8+5+4+9}{5} = \frac{30}{5} = 6.$$

New centre = (6, 6)

$$C_3 = \left(\frac{2+1}{2}, \frac{5+2}{2} \right) = (1.5, 3.5)$$

Ans:-

After 1 epoch, centres
are

C₁ : (2, 10)

C₂ : (6, 6)

C₃ : (1.5, 3.5)



DATE

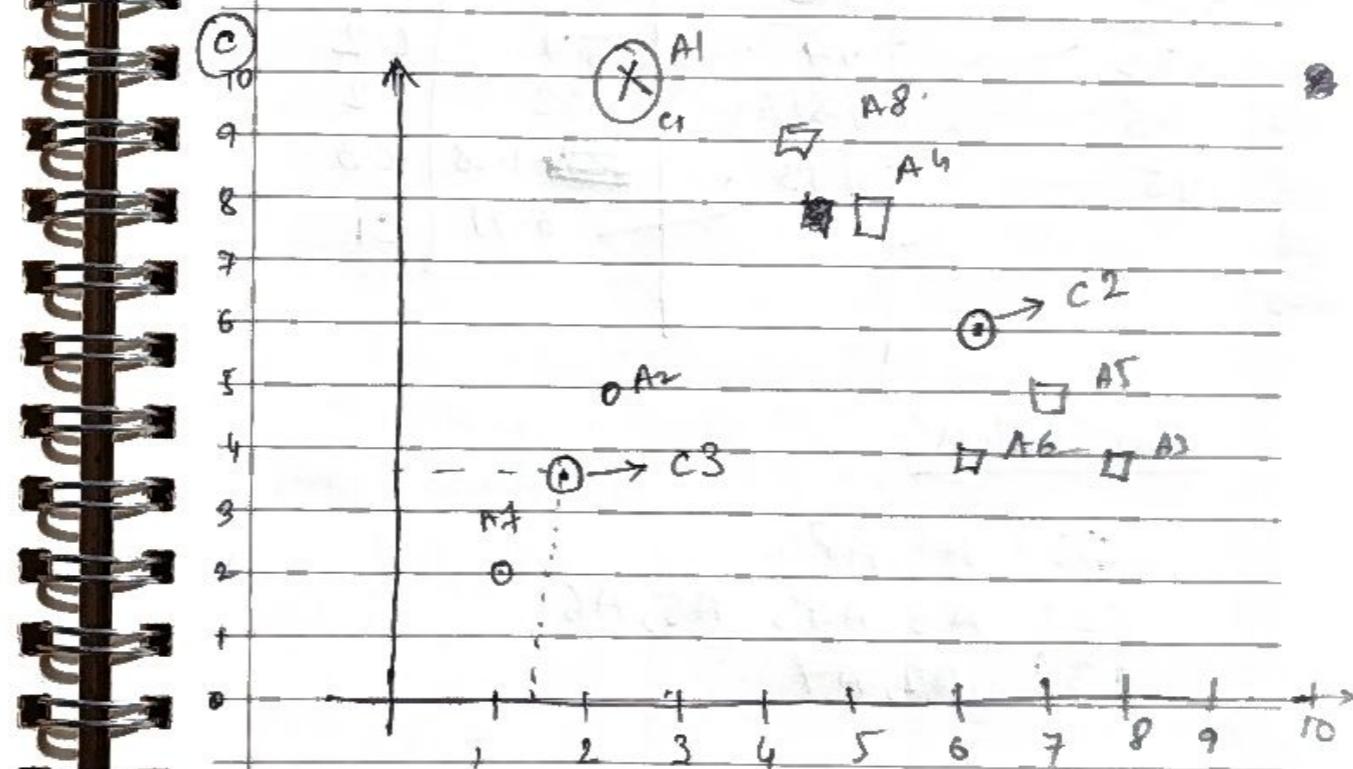
PLACE

and clusters are

(A₁), (A₃, A₄, A₅, A₆, A₈)

and (A₂, A₇).

X → cluster 1 O → cluster 3
□ → cluster 2



C₁, C₂, C₃ are the cluster centres.



DATE

PLACE

Epoch 2

PT	to C1	DIST to C2	DIST to C3	Assigned cluster
				C1
A1	0			
A2	$\sqrt{23} = 5$	$\sqrt{16+16} = \sqrt{32}$	$\sqrt{42.25} = 6.5$	C1
A3	$\sqrt{72} = 8.49$	$\sqrt{8}$	1.58	C3
A4	$\sqrt{13} = 3.61$	$\sqrt{5}$	6.52	C2
A5	$\sqrt{50}$	$\sqrt{20}$	5.7	C2
A6	$\sqrt{52}$	$\sqrt{4}$	5.7	C2
A7	$\sqrt{45}$	$\sqrt{28+16} = 6.4$	4.53	C2
A8	$\sqrt{5}$	$\sqrt{13}$	1.58 6.71	C3

New Clusters:

C1: A1, A8

C2: A3, A4, A5, A6.

C3: A2, A7.

New centre means:-

$$C1 \Rightarrow \frac{2+4}{2}, \frac{10+9}{2} = (3, 9.5)$$

$$C2 \Rightarrow \frac{8+5+7+6}{4}, \frac{4+8+5+4}{4} = (6.5, 5.25)$$

$$C3 \Rightarrow (1.5, 3.5)$$



DATE

PLACE

Compute distances again.

PT	C1	C2	C3	cluster
A1	1.80	6.36	6.5	C1
A2	4.6	4.76	1.58	C3
A3	6.72	2.06	~6.52	C2
A4	2.5	3.04	~5.7	C1 & changed
A5	4.84	1.45	~5.7	C2
A6	6.02	1.53	~4.53	C2
A7	8.83	6.2	1.58	C3
A8	1.12	4.2	6.7	C1

New clusters: C1 \rightarrow A1, A4, A8.

C2: A3, A5, A6.

C3: A2, A7,

New centres:

$$C1 = \frac{2+5+4}{3}, \frac{10+8+9}{3} = (3.67, 9)$$

$$C2 = \left(\frac{8+5+7+6}{4}, \frac{4+8+5+4}{4} \right) = (6.5, 5.25)$$

$$C3 = (1.5, 3.5) \text{ same}$$

Thus, convergence is reached after epoch 4.



DATE _____

PLACE _____

final cluster Results

Epoch	cluster 1	cluster 2	Cluster 3
1	A1	A3, A4, A5, A6, A8	A2, A7
2	A1, A8	A3, A4, A5, A6	A2, A7
3	A1, A6, A8	A3, A5, A6	A2, A7
4	A1, A4, A8	A3, A5, A6	A2, A7



DATE _____

PLACE _____

(e) we need to prove

$$\sum_{i=1}^m \|z_i - z\|^2 \geq \sum_{i=1}^m \|z_i - \bar{z}\|^2$$

This is equivalent to saying mean \bar{z} minimizes the sum of squared distances from points z_i . Any other point gives greater or equal distance squared.

we know,

$$\text{mean } \bar{z} = \frac{1}{m} \sum_{i=1}^m z_i$$

 z = arbitrary number

using an identity,

$$\begin{aligned} \|a - b\|^2 &= \|a - c + c - b\|^2 \\ &\geq \|a - c\|^2 + 2 \langle a - c, c - b \rangle \\ &\quad + \|c - b\|^2 \end{aligned}$$



DATE _____

PLACE _____

$$\text{When, } a = z_i$$

$$b = \bar{z}$$

$$c = \bar{z}$$

$$\sum_{i=1}^m \|z_i - z\|^2 = \sum_{i=1}^m \|z_i - \bar{z}\|^2 + 2 \sum_{i=1}^m \langle z_i - \bar{z}, \bar{z} - z \rangle$$

$$+ \sum_{i=1}^m \|\bar{z} - z\|^2.$$

Now,

$$\sum_{i=1}^m \langle z_i - \bar{z}, \bar{z} - z \rangle$$

$$= \sum_{i=1}^m \langle (z_i - \bar{z}), (\bar{z} - z) \rangle$$

$$= \left\langle \sum_{i=1}^m (z_i - \bar{z}), \bar{z} - z \right\rangle$$

$$\underbrace{\quad}_{\downarrow} \quad 0$$

$$= \langle 0, \bar{z} - z \rangle = 0.$$

Mean balances out the data pts. So sum of all offsets is always 0.



DATE _____

PLACE _____

$$\sum_{i=1}^m \|\bar{z} - z\|^2 = m \cdot \|\bar{z} - z\|^2$$

Thus, expression boils down to,

$$\begin{aligned} \sum_{i=1}^m \|z_i - z\|^2 &= \sum_{i=1}^m \|z_i - \bar{z}\|^2 \\ &\quad + m \cdot \|\bar{z} - z\|^2 \\ &\Rightarrow \sum_{i=1}^m \|z_i - \bar{z}\|^2 \end{aligned}$$

with equality m when $z = \bar{z}$

The second term is always non-negative since, its squared norm:



DATE _____

PLACE _____

- (f) The process is a description of K-means with a ~~as~~ greedy initialization.

1st centre: random

2nd centre: pick point farthest from existing centers.

After that standard K-means.

At each step,

Given centroids μ^t , each pt is assigned to nearest centroid.

Here, sum of square errors doesn't increase as each pt moves closer to its actual centre.

In update step,

Given C^{t+1} clusters, we recompute new centroids μ^{t+1} as the mean of each cluster.

Earlier we proved, mean minimises sum of squared distances.



DATE _____

PLACE _____

Thus, SSE decreases or remains same.

$$\text{SSE}(C^{t+1}, \mu^{t+1}) \leq \text{SSE}(C^t, \mu^t).$$

thus, as long as the clustering changes, we will move towards convergence with greedy initialization.

Greedy initialization tries to spread out centres in the start preventing any kind of skewed choice.

This results in a lower iterations need to converge, since we have reduced chances of local bad minima, when compared to random initialization.

- (c) Each point can belong to any one of the K-clusters.

Thus, K^n possible ways to assign n points to K clusters.

Thus, each iteration has either to decrease SSE strictly(until convergence).



DATE _____

PLACE _____

Since, there are finitely many clusterings, it has to terminate in finite loop. Thus, guarantee termination.

At first initialisation phase,

clustering starts at random locations.

events become of high pre-

dictability and so often will

occur some small separation between

clusters and go through random

relocation rather than growth.

One pair of particles may start

as islands and

consequently of low sticking rate and

islands will

not cluster and eventually just

remain as individual clusters.



DATE 11/04/25

PLACE

- (4) Given, $k = 2$ for a dataset with 4 points
 (5) & GMM to be fit.

at t^{th} time step, we have :

θ_t with following parameters:

$$\pi_1 = 0.3 \quad \pi_2 = 0.7$$

$$\mu_1 = 2 \quad \mu_2 = 3$$

$$\sigma_1^2 = 1 \quad \sigma_2^2 = 1$$

$$f(x_i | z_i = 1) \text{ & } f(x_i | z_i = 2)$$

for $i = 1 \text{ to } 4$ are given.

In E step, we are calculating responsibilities γ_i^k , the probability that point x_i belongs to cluster K .

$z_i \rightarrow$ represents the latent variable that indicates which ~~cluster~~ Gaussian component generated the data point x_i



DATE

PLACE

Responsibility of cluster K for point i is

$$\gamma_i^K = \frac{\pi_k \cdot f(x_i | z_i = k)}{\sum_{j=1}^K \pi_j \cdot f(x_i | z_i = j)}$$

Responsibility of cluster 2 from point 1 ie. how likely is that point 1 comes from cluster 2 is

$$i = 1, K = 2$$

$$(0.7) (0.054)$$

$$\gamma_1^2 = \frac{0.3 \times 0.242}{0.3 \times 0.242 + 0.7 \times 0.054}$$

$$0.0378$$

$$= \frac{0.0378}{0.0726 + 0.0378} = 0.3423$$

So, probability of point 1 belonging to cluster 1 is 34.23%.



DATE _____
PLACE _____

(b) if we pause after cluster E step, we have all the responsibilities π_i^k .

If we do a hard assignment, we simply shall assign each point to the cluster for which it has highest responsibility.

i.e. $z_i = \arg \max_k \pi_i^k$

Since, $\pi_1^2 = 0.3423$

$$\pi_1^1 = 1 - 0.3423 = \frac{0.6577}{0.6577}$$

Among this, max is cluster 1.

So, $\boxed{z_1 = 1}$



DATE _____
PLACE _____

Point 2

cluster 1 point 2

$$\pi_2^1 = \frac{\pi_1 f(x_2 | z_2 = 1)}{\sum_{j=1}^K \pi_j f(x_2 | z_2 = j)}$$

$$= \frac{\pi_1 f(x_2 | z_2 = 1)}{\pi_1 f(x_2 | z_2 = 1) + \pi_2 f(x_2 | z_2 = 2)}$$

$$= \frac{0.3 \times 0.399}{0.3 \times 0.399 + 0.7 \times 0.242}$$

$$= \frac{0.1197}{0.2891} = 0.4141$$

Here, $\pi_2^2 = 1 - 0.4141 > \pi_2^1$

thus $\boxed{z_2 = 2}$.



DATE _____

PLACE _____

$$\lambda_3^1 = \frac{\pi_1 f(x_3 | z_3=1)}{\pi_1 f(x_3 | z_3=1) + \pi_2 f(x_3 | z_3=2)}$$

$$= \frac{0.3 \times 0.242}{0.3 \times 0.242 + 0.7 \times 0.399}$$

$$= \frac{0.0726}{0.3519} = 0.2063.$$

$$\lambda_3^2 = 1 - 0.2063 = 0.7937$$

Thus, ~~etc~~ z₃ = 2

Point 4

$$\lambda_4^1 = \frac{\pi_1 f(x_4 | z_4=1)}{\pi_1 f(x_4 | z_4=1) + \pi_2 f(x_4 | z_4=2)}$$

$$= \frac{0.3 \times 0.054}{0.3 \times 0.054 + 0.7 \times 0.242}$$

$$= \frac{0.0162}{0.1856} = 0.0873$$



DATE _____

PLACE _____

$$\lambda_4^2 = 1 - \lambda_4^1 \geq \lambda_4^1$$

Since $\lambda_4^1 = 0.0873$

thus $z_4 = 2$.

Final vector notation of cluster

$$= [1 \ 2 \ 2 \ 2]$$